

Introduction to Data Visualization in Python

 towardsdatascience.com/introduction-to-data-visualization-in-python-89a54c97fbed

March 6, 2019

How to make graphs using Matplotlib, Pandas and Seaborn



Figure 1: Photo by on

Data visualization is the discipline of trying to understand data by placing it in a visual context so that patterns, trends and correlations that might not otherwise be detected can be exposed.

Python offers multiple great graphing libraries that come packed with lots of different features. No matter if you want to create interactive, live or highly customized plots python has an excellent library for you.

To get a little overview here are a few popular plotting libraries:

- low level, provides lots of freedom
- easy to use interface, built on Matplotlib
- high-level interface, great default styles
- based on R's ggplot2, uses
- can create interactive plots

In this article, we will learn how to create basic plots using Matplotlib, Pandas visualization and Seaborn as well as how to use some specific features of each library. This article will focus on the syntax and not on interpreting the graphs, which I will cover in another blog post.

In further articles, I will go over interactive plotting tools like Plotly, which is built on D3 and can also be used with JavaScript.

Importing Datasets

In this article, we will use two datasets which are freely available. The Iris and Wine Reviews dataset, which we can both load in using pandas `read_csv` method.

```
iris = pd.read_csv('iris.csv')
```

	sepal_length	sepal_width	petal_length	petal_width	class
0	5.1	3.5	1.4	0.2	Iris-setosa
1	4.9	3.0	1.4	0.2	Iris-setosa
2	4.7	3.2	1.3	0.2	Iris-setosa
3	4.6	3.1	1.5	0.2	Iris-setosa
4	5.0	3.6	1.4	0.2	Iris-setosa

Figure 2: Iris dataset head

	country	description	designation	points	price	province	region_1	region_2	taster_name	taster_twitter_handle	title	variety	winery
0	Italy	Aromas include tropical fruit, broom, brimston...	Vulkà Bianco	87	NaN	Sicily & Sardinia	Etna	NaN	Kerin O'Keefe	@kerinokeefe	Nicosia 2013 Vulkà Bianco (Etna)	White Blend	Nicosia
1	Portugal	This is ripe and fruity, a wine that is smooth...	Avidagos	87	15.0	Douro	NaN	NaN	Roger Voss	@vossroger	Quinta dos Avidagos 2011 Avidagos Red (Douro)	Portuguese Red	Quinta dos Avidagos
2	US	Tart and snappy, the flavors of lime flesh and...	NaN	87	14.0	Oregon	Willamette Valley	Willamette Valley	Paul Gregutt	@paulgwine	Rainstorm 2013 Pinot Gris (Willamette Valley)	Pinot Gris	Rainstorm
3	US	Pineapple rind, lemon pith and orange blossom ...	Reserve Late Harvest	87	13.0	Michigan	Lake Michigan Shore	NaN	Alexander Peartree	NaN	St. Julian 2013 Reserve Late Harvest Riesling ...	Riesling	St. Julian
4	US	Much like the regular bottling from 2012, this...	Vintner's Reserve Wild Child Block	87	65.0	Oregon	Willamette Valley	Willamette Valley	Paul Gregutt	@paulgwine	Sweet Cheeks 2012 Vintner's Reserve Wild Child...	Pinot Noir	Sweet Cheeks

Figure 3: Wine Review dataset head

Matplotlib

Matplotlib is the most popular python plotting library. It is a low-level library with a Matlab like interface which offers lots of freedom at the cost of having to write more code.

To install Matplotlib pip and conda can be used.

```
pip install matplotlib
or
conda install matplotlib
```

Matplotlib is specifically good for creating basic graphs like line charts, bar charts, histograms and many more. It can be imported by typing:

```
import matplotlib.pyplot as plt
```

Scatter Plot

To create a scatter plot in Matplotlib we can use the `scatter` method. We will also create a figure and an axis using `plt.subplots` so we can give our plot a title and labels.

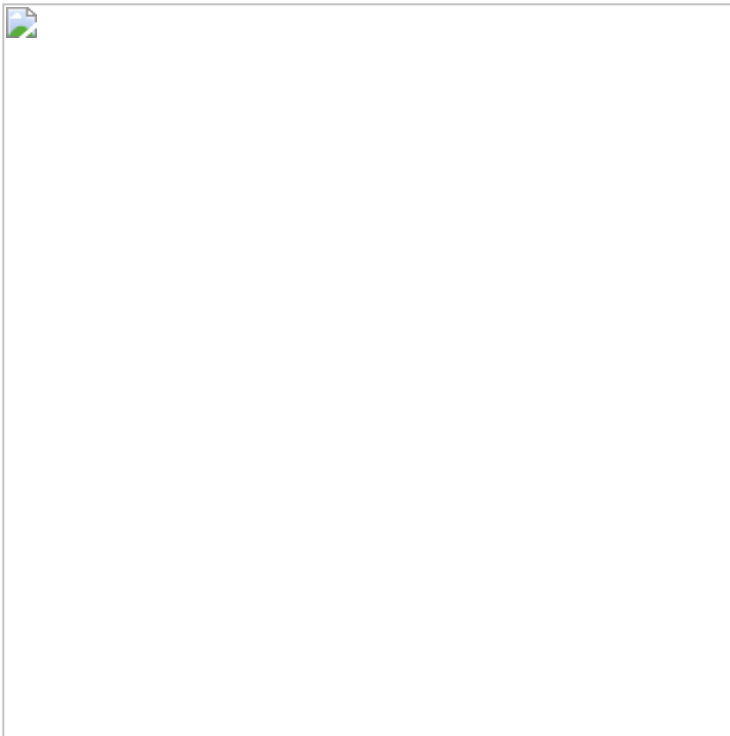
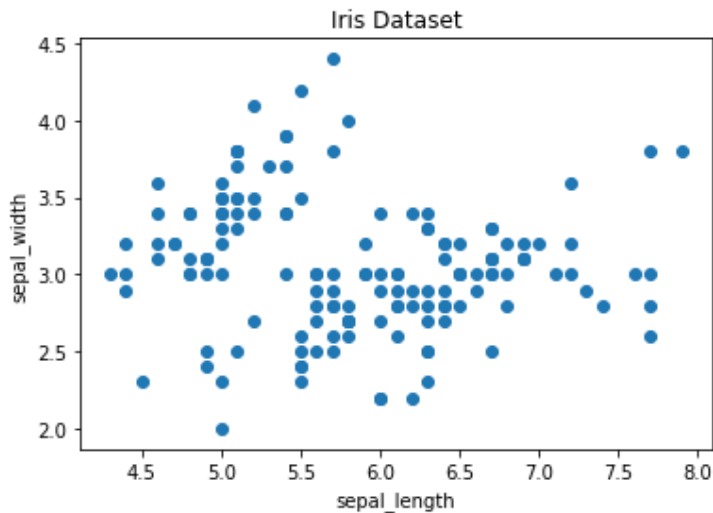


Figure 4: Matplotlib Scatter plot

We can give the graph more meaning by coloring in each data-point by its class. This can be done by creating a dictionary which maps from class to color and then scattering each point on its own using a for-loop and passing the respective color.

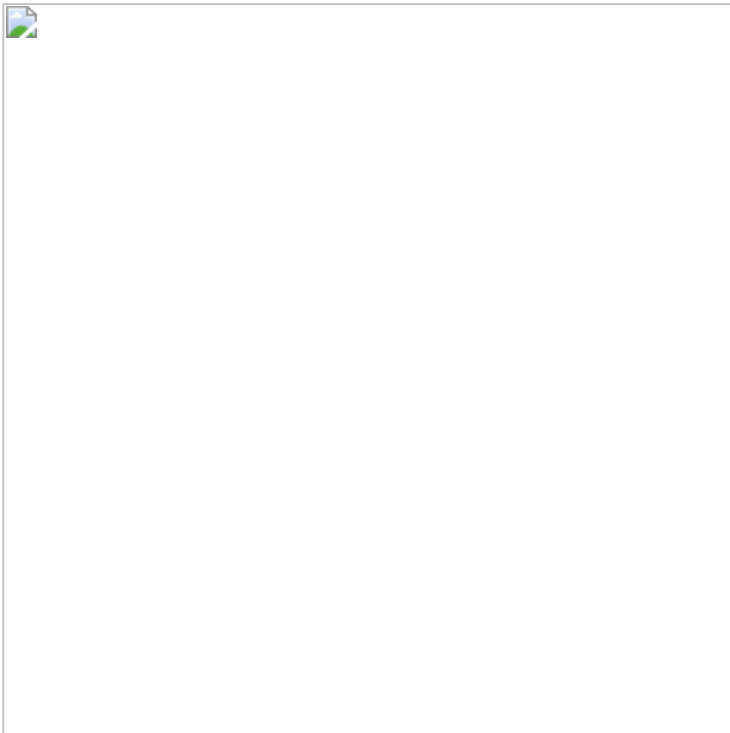
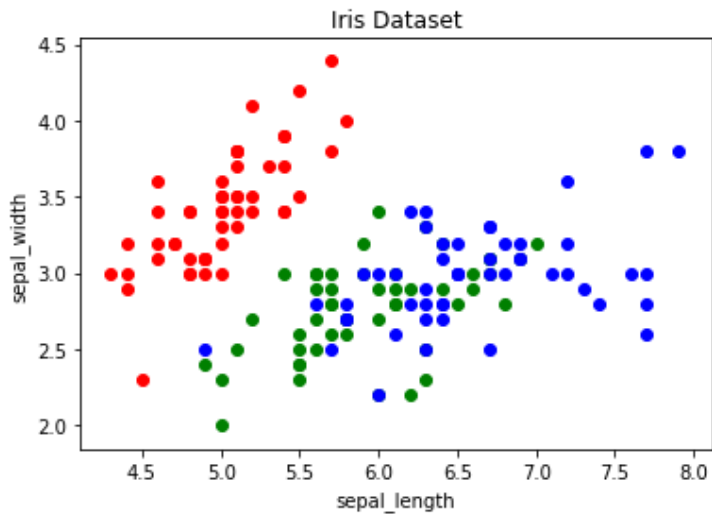


Figure 5: Scatter Plot colored by class

Line Chart

In Matplotlib we can create a line chart by calling the `plot` method. We can also plot multiple columns in one graph, by looping through the columns we want and plotting each column on the same axis.

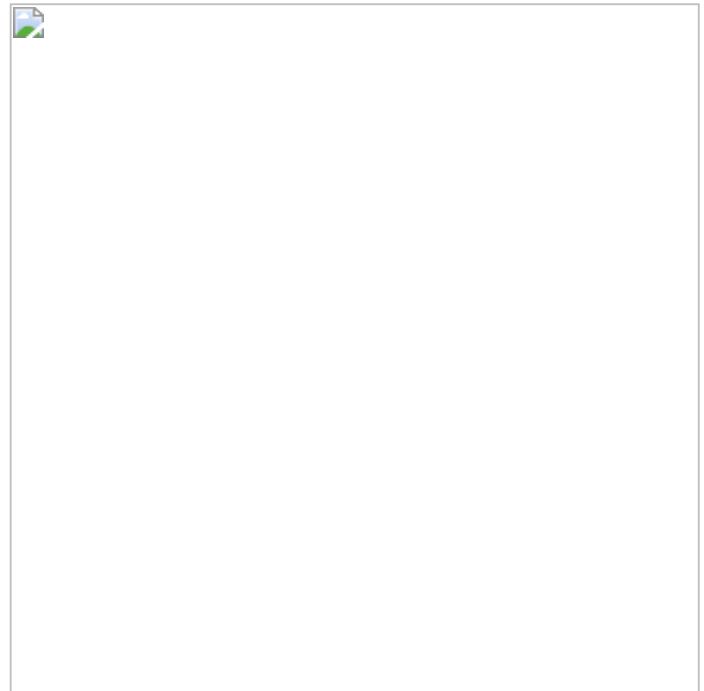
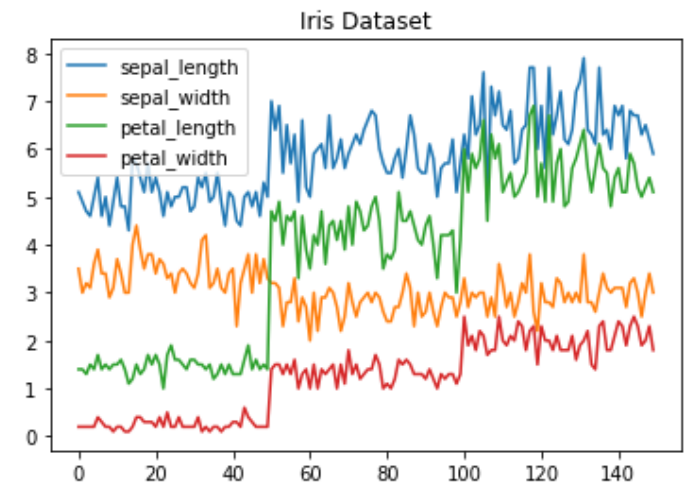


Figure 6: Line Chart

Histogram

In Matplotlib we can create a Histogram using the `hist` method. If we pass it categorical data like the points column from the wine-review dataset it will automatically calculate how often each class occurs.

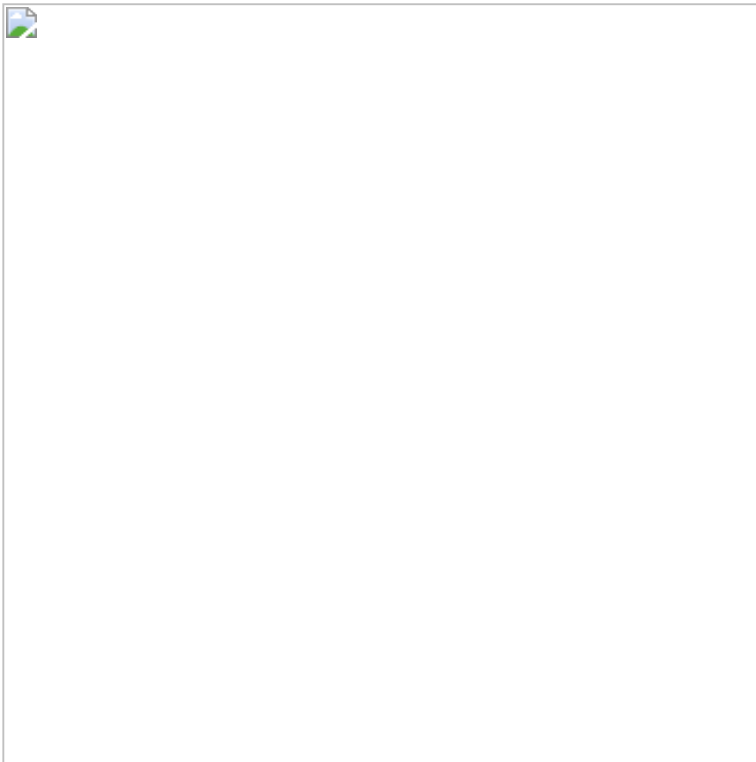
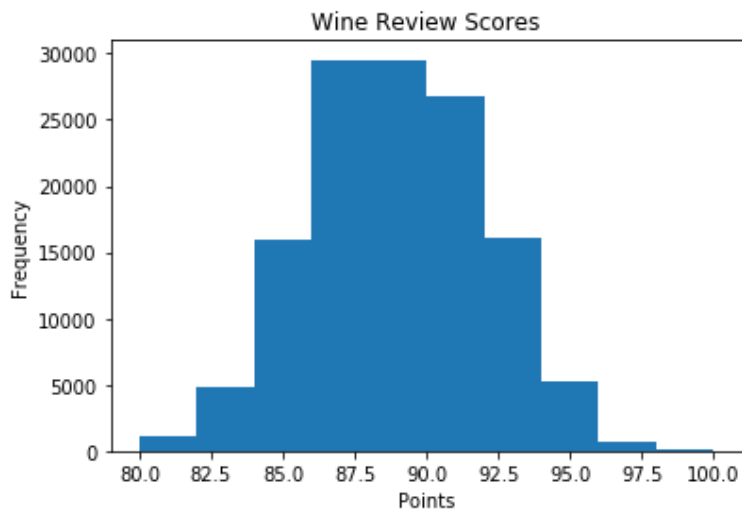


Figure 7: Histogram

Bar Chart

A bar chart can be created using the `bar` method. The bar-chart isn't automatically calculating the frequency of a category so we are going to use pandas `value_counts` function to do this. The bar-chart is useful for categorical data that doesn't have a lot of different categories (less than 30) because else it can get quite messy.

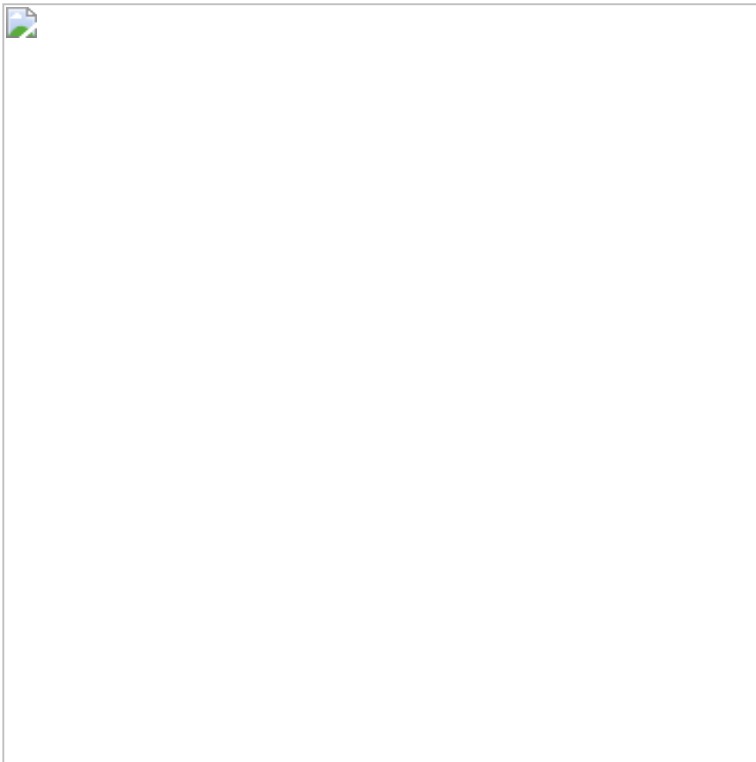
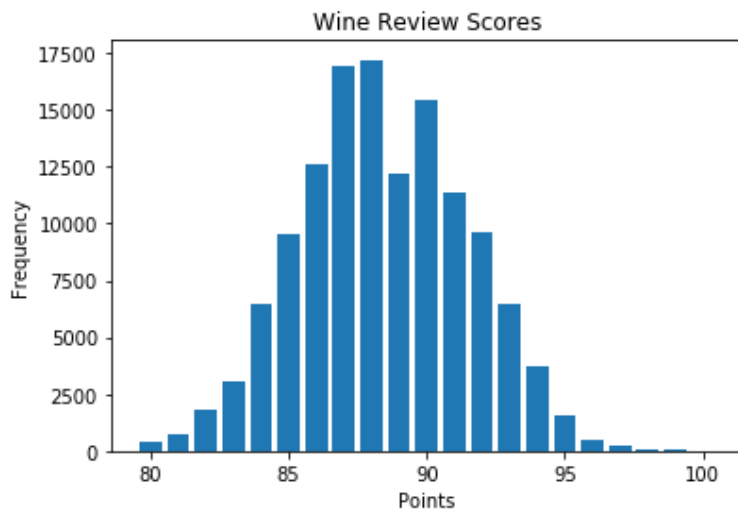


Figure 8: Bar-Chart

Pandas Visualization

Pandas is an open source high-performance, easy-to-use library providing data structures, such as dataframes, and data analysis tools like the visualization tools we will use in this article.

Pandas Visualization makes it really easy to create plots out of a pandas dataframe and series. It also has a higher level API than Matplotlib and therefore we need less code for the same results.

Pandas can be installed using either pip or conda.

```
pip install pandas
or
conda install pandas
```

Scatter Plot

To create a scatter plot in Pandas we can call `<dataset>.plot.scatter()` and pass it two arguments, the name of the x-column as well as the name of the y-column. Optionally we can also pass it a title.

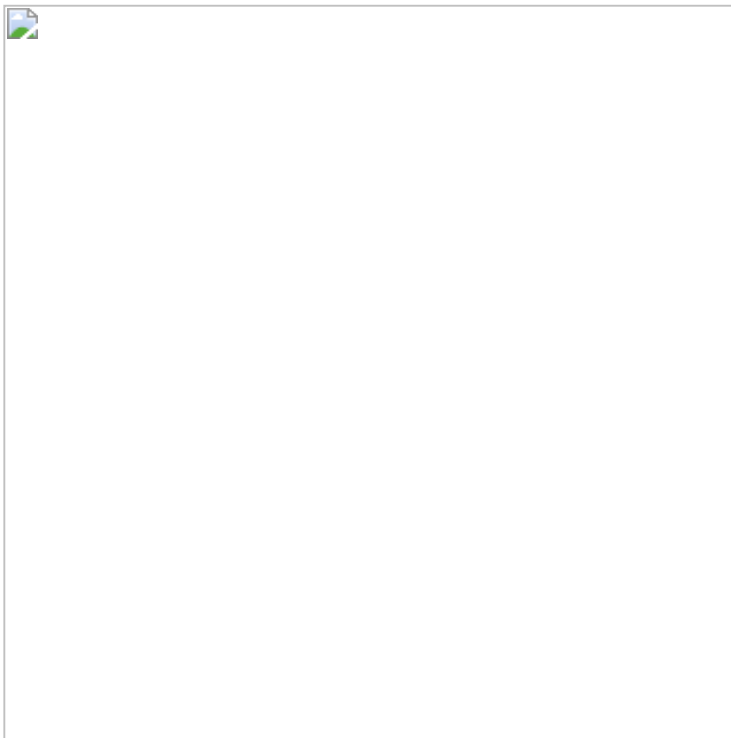
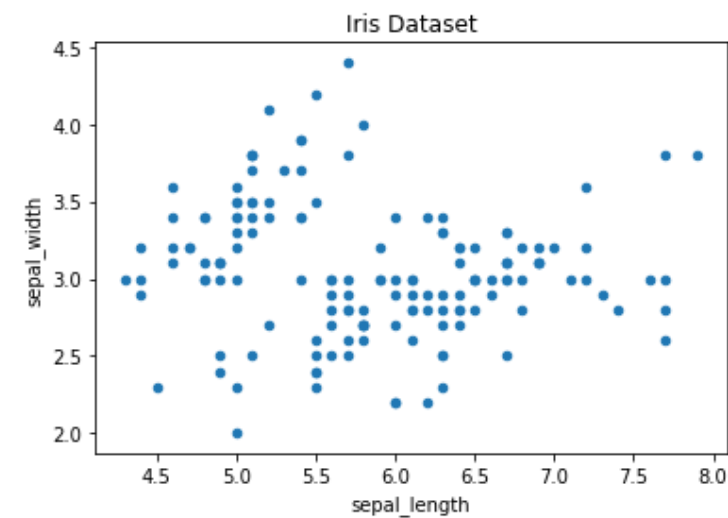


Figure 9: Scatter Plot

As you can see in the image it is automatically setting the x and y label to the column names.

Line Chart

To create a line-chart in Pandas we can call `<dataframe>.plot.line()`. Whilst in Matplotlib we needed to loop-through each column we wanted to plot, in Pandas we don't need to do this because it automatically plots all available numeric columns (at least if we don't specify a specific column/s).

If we have more than one feature Pandas automatically creates a legend for us, as can be seen in the image above.



Figure 10: Line Chart

Histogram

In Pandas, we can create a Histogram with the `plot.hist` method. There aren't any required arguments but we can optionally pass some like the bin size.

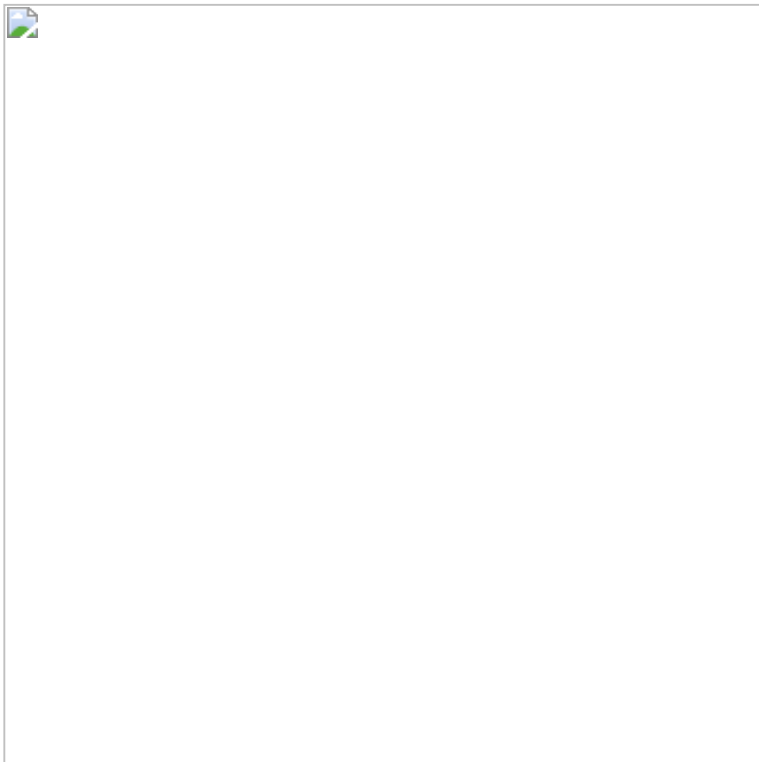
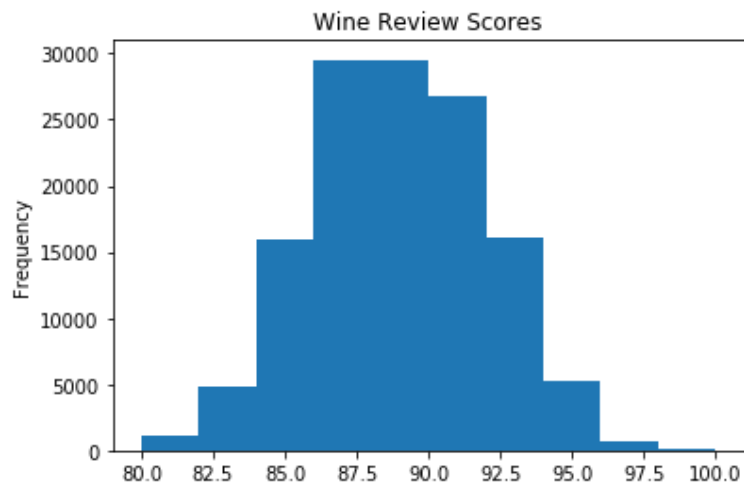


Figure 11: Histogram

It's also really easy to create multiple histograms.

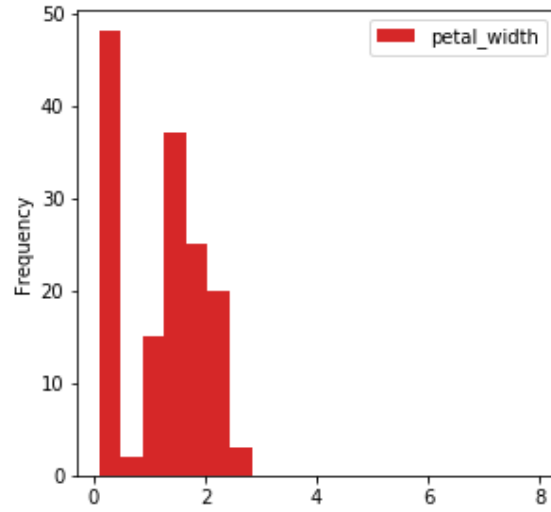
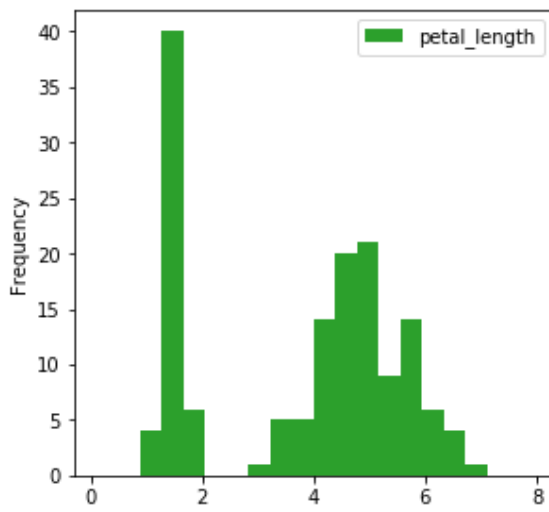
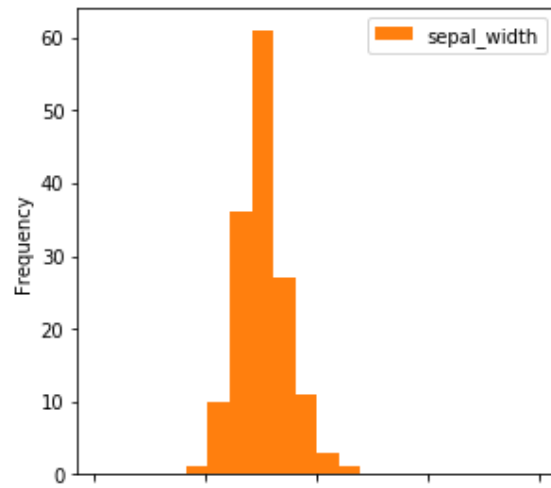
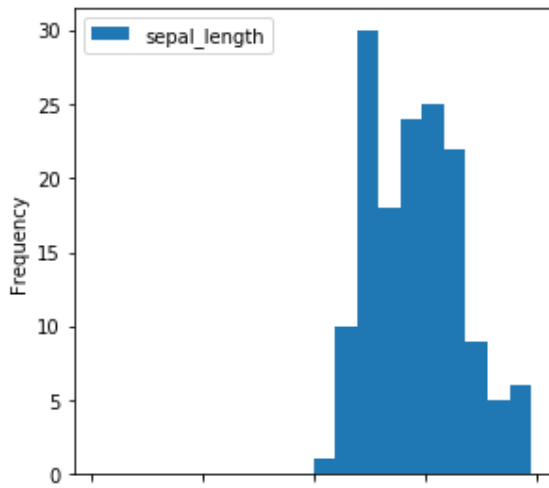




Figure 12: Multiple Histograms

The `subplots` argument specifies that we want a separate plot for each feature and the `layout` specifies the number of plots per row and column.

Bar Chart

To plot a bar-chart we can use the `plot.bar()` method, but before we can call this we need to get our data. For this we will first count the occurrences using the `value_count()` method and then sort the occurrences from smallest to largest using the `sort_index()` method.

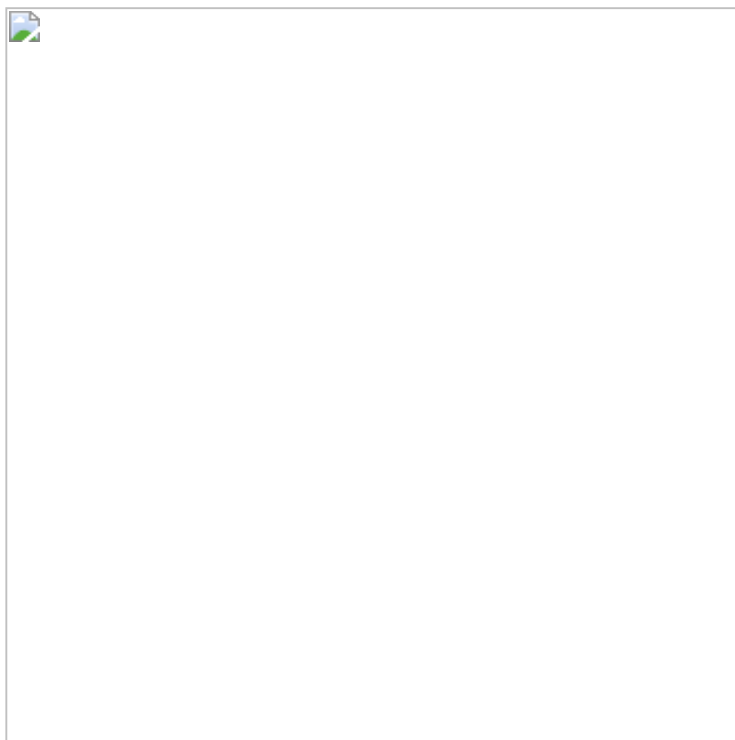
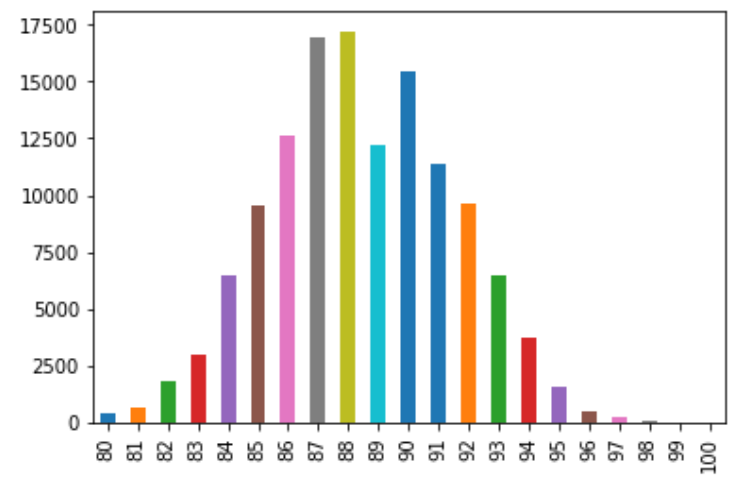


Figure 13: Vertical Bar-Chart

It's also really simple to make a horizontal bar-chart using the `plot.barh()` method.

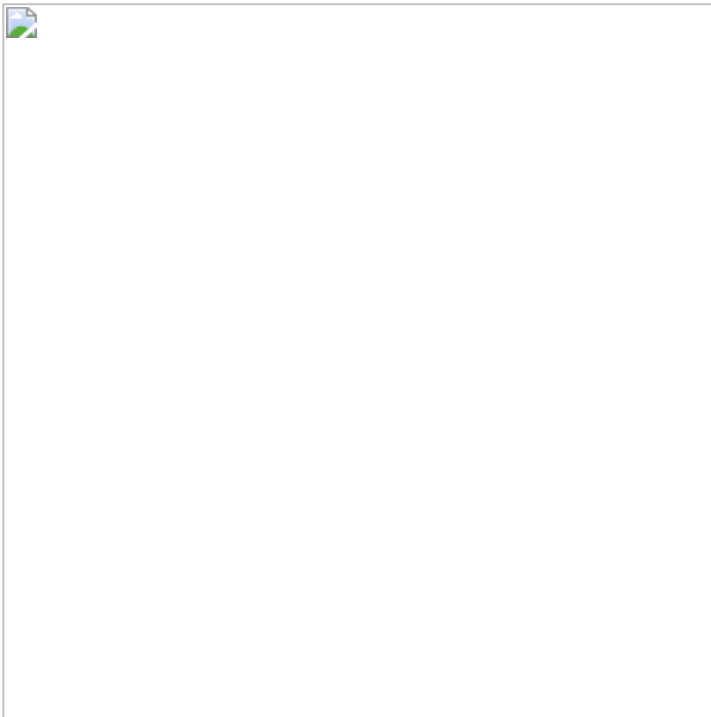
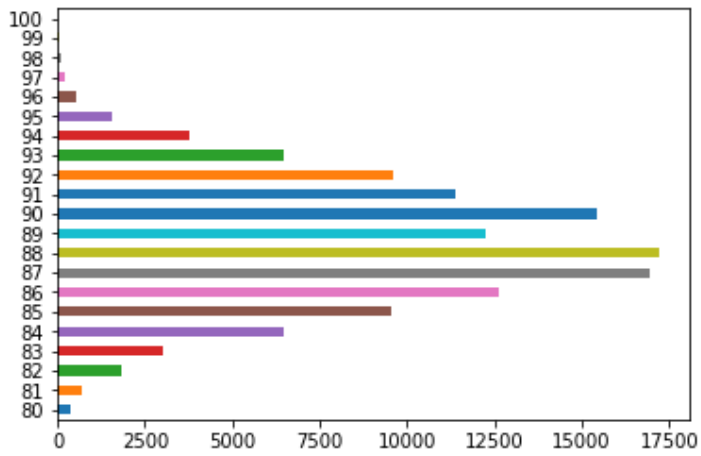


Figure 14: Horizontal Bar-Chart

We can also plot other data then the number of occurrences.

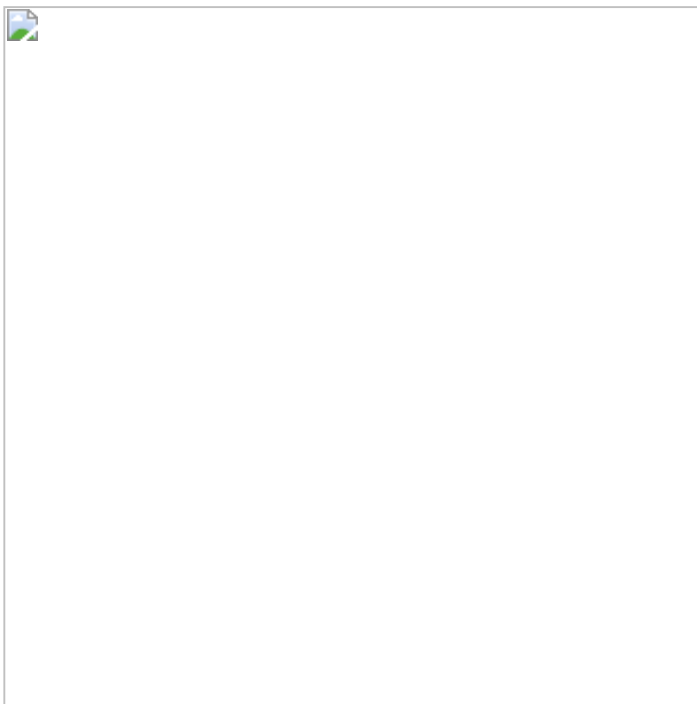
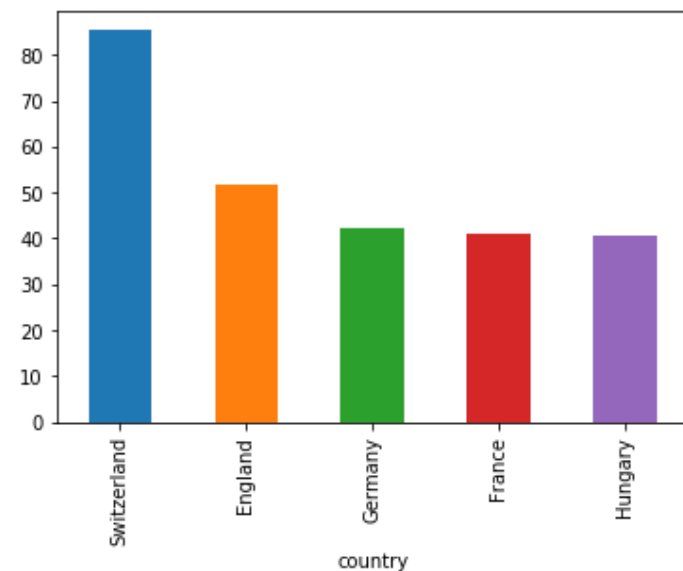


Figure 15: Countries with the most expensive wine(on average)

In the example above we grouped the data by country and then took the mean of the wine prices, ordered it, and plotted the 5 countries with the highest average wine price.

Seaborn

Seaborn is a Python data visualization library based on Matplotlib. It provides a high-level interface for creating attractive graphs.

Seaborn has a lot to offer. You can create graphs in one line that would take you multiple tens of lines in Matplotlib. Its standard designs are awesome and it also has a nice interface for working with pandas dataframes.

It can be imported by typing:

```
import seaborn as sns
```

Scatter plot

We can use the `.scatterplot` method for creating a scatterplot, and just as in Pandas we need to pass it the column names of the x and y data, but now we also need to pass the data as an additional argument because we aren't calling the function on the data directly as we did in Pandas.

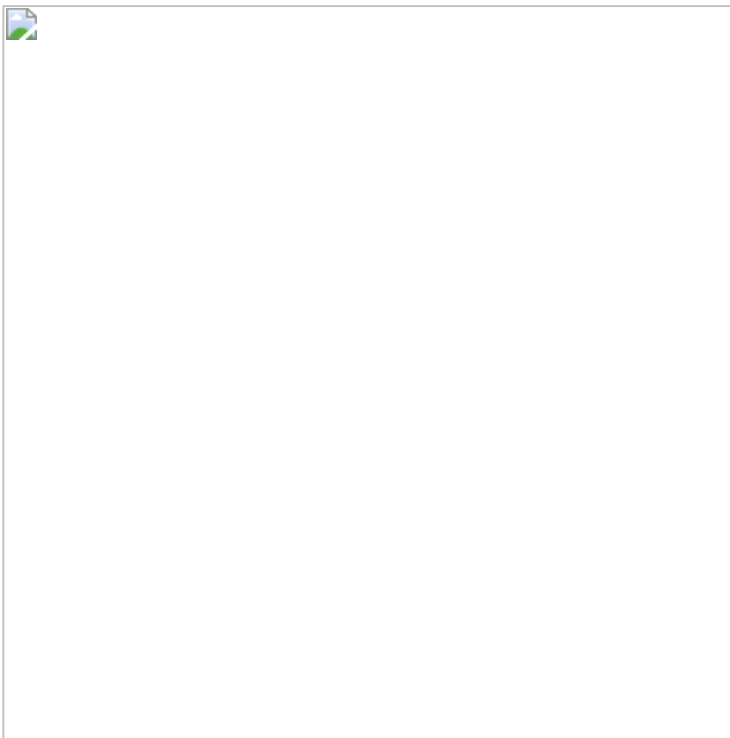
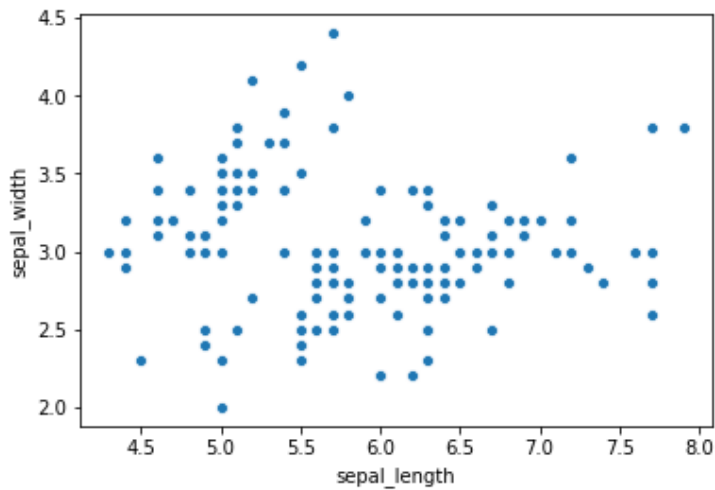


Figure 16: Scatterplot

We can also highlight the points by class using the `hue` argument, which is a lot easier than in Matplotlib.

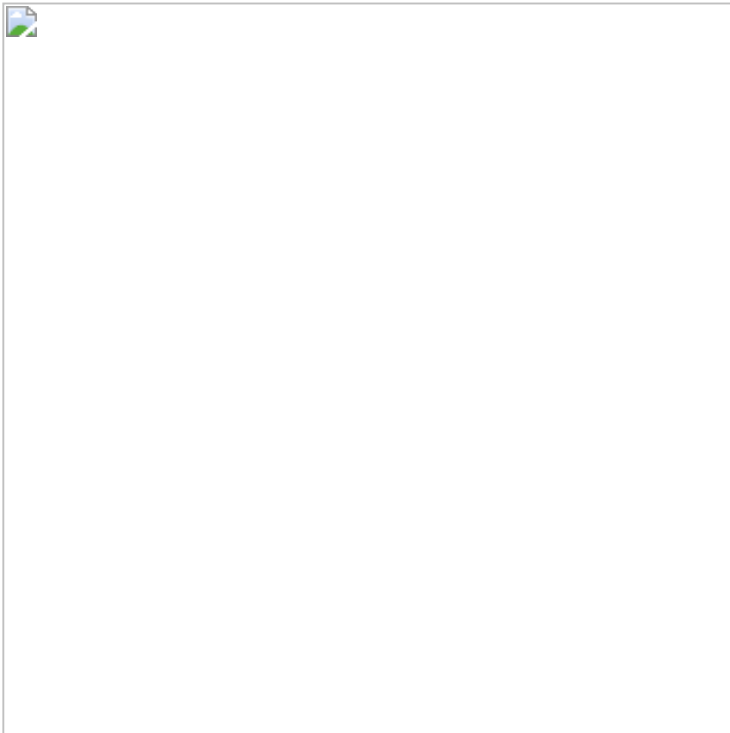
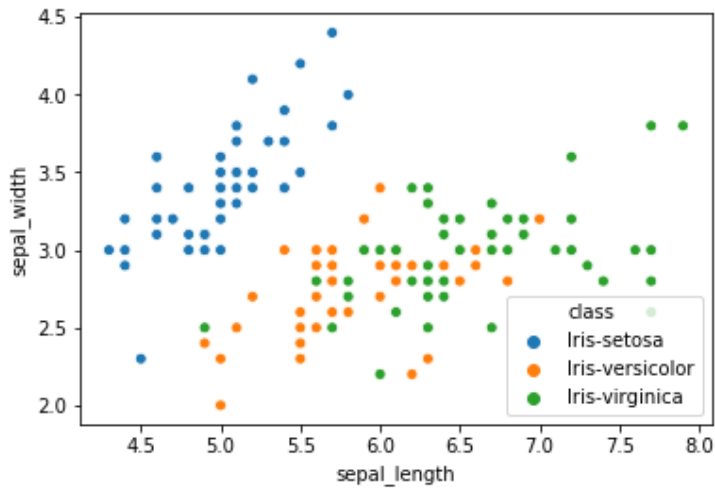


Figure 17: Scatterplot colored by class

Line chart

To create a line-chart the `sns.lineplot` method can be used. The only required argument is the data, which in our case are the four numeric columns from the Iris dataset. We could also use the `sns.kdeplot` method which rounds of the edges of the curves and therefore is cleaner if you have a lot of outliers in your dataset.

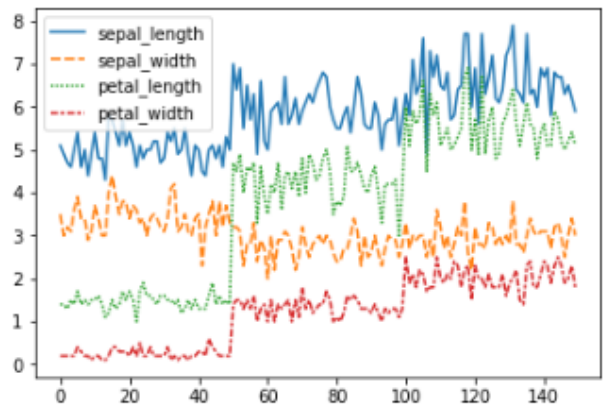


Figure 18: Line Chart

Histogram

To create a histogram in Seaborn we use the `sns.distplot` method. We need to pass it the column we want to plot and it will calculate the occurrences itself. We can also pass it the number of bins, and if we want to plot a gaussian kernel density estimate inside the graph.

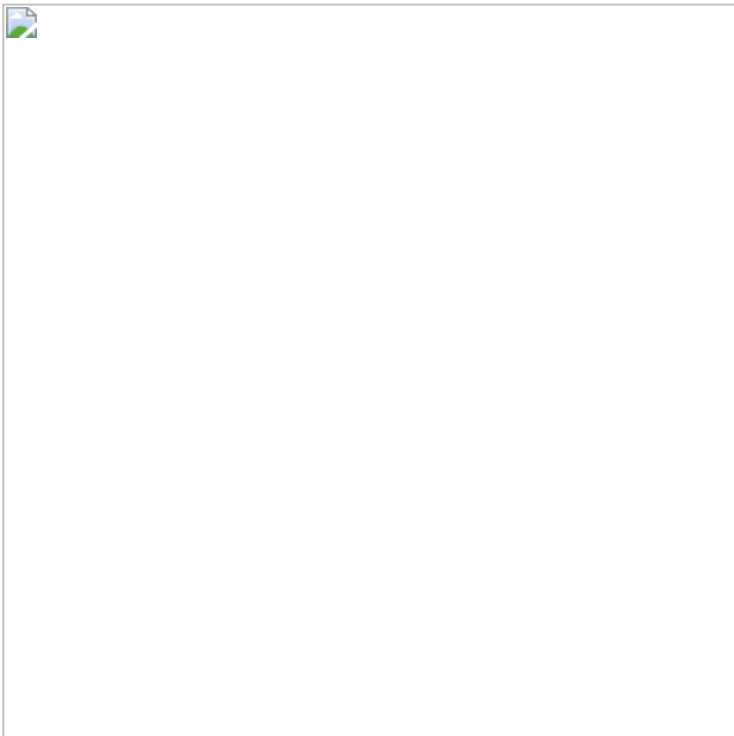
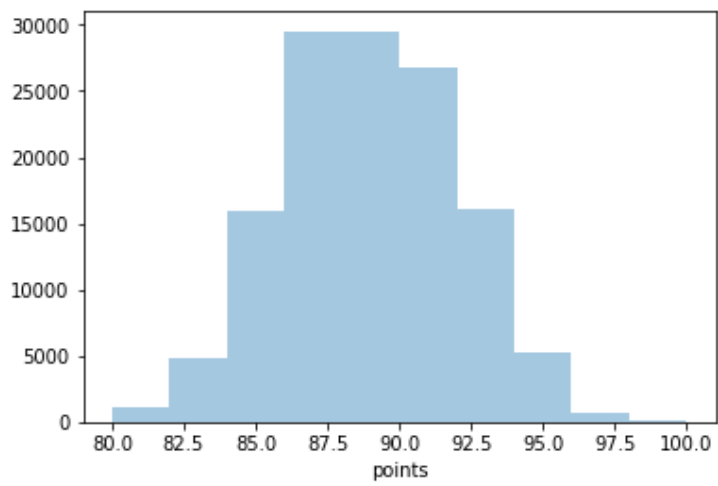


Figure 19: Histogram

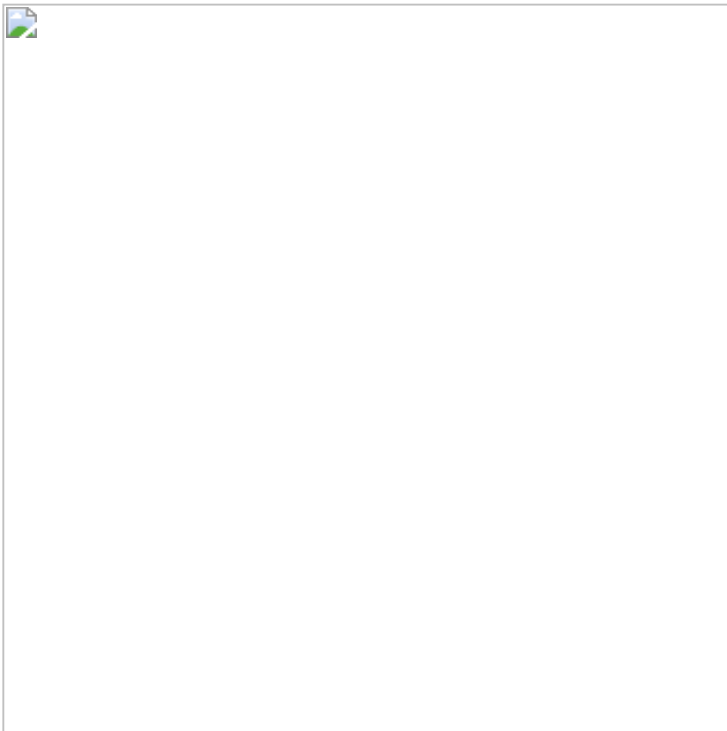
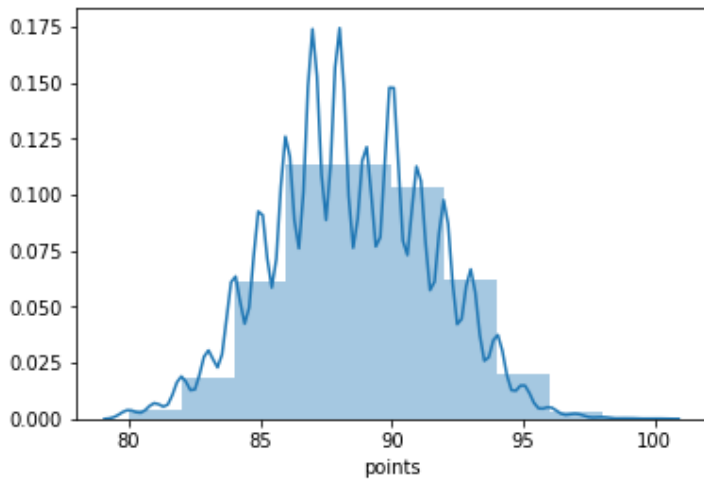


Figure 20: Histogram with gaussian kernel density estimate

Bar chart

In Seaborn a bar-chart can be created using the `sns.countplot` method and passing it the data.

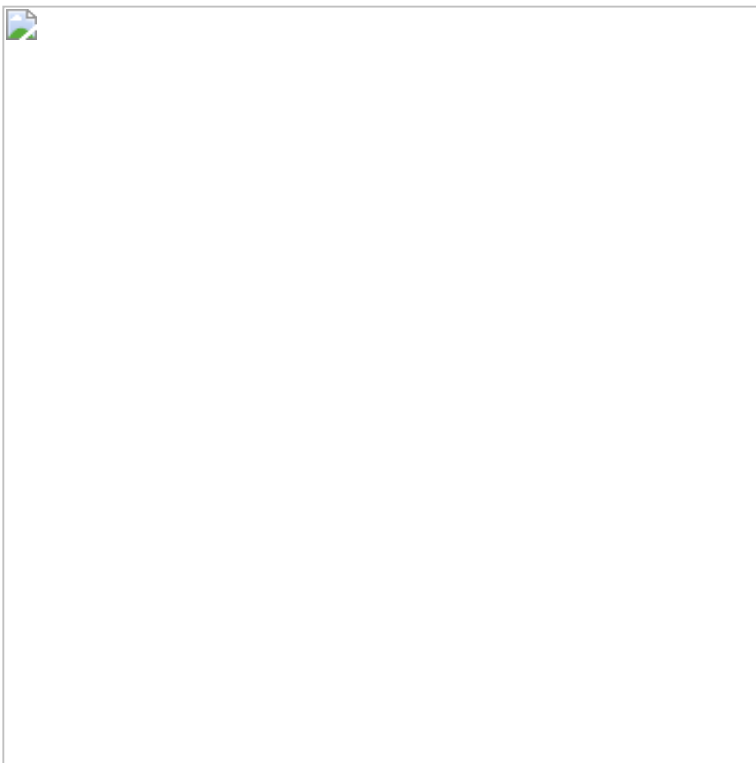
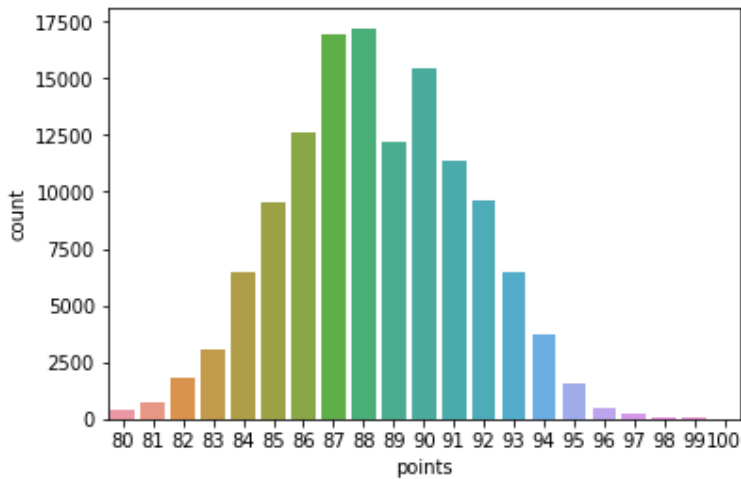


Figure 21: Bar-Chart

Other graphs

Now that you have a basic understanding of the Matplotlib, Pandas Visualization and Seaborn syntax I want to show you a few other graph types that are useful for extracting insides.

For most of them, Seaborn is the go-to library because of its high-level interface that allows for the creation of beautiful graphs in just a few lines of code.

Box plots

A Box Plot is a graphical method of displaying the five-number summary. We can create box plots using seaborn's `sns.boxplot` method and passing it the data as well as the x and y column name.

Box Plots, just like bar-charts are great for data with only a few categories but can get messy really quickly.

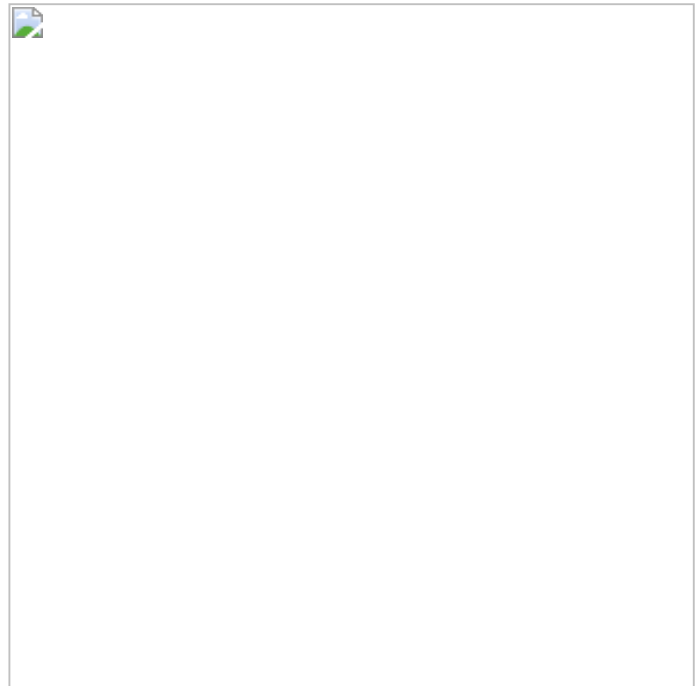
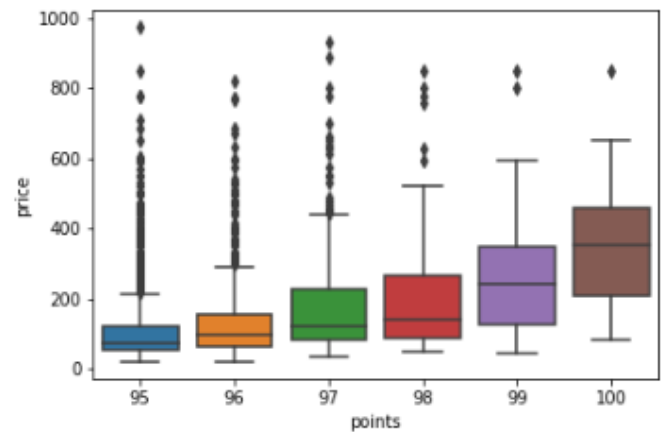


Figure 22: Boxplot

Heatmap

A Heatmap is a graphical representation of data where the individual values contained in a matrix are represented as colors. Heatmaps are perfect for exploring the correlation of features in a dataset.

To get the correlation of the features inside a dataset we can call `<dataset>.corr()`, which is a Pandas dataframe method. This will give us the correlation matrix.

We can now use either Matplotlib or Seaborn to create the heatmap.

Matplotlib:

To add annotations to the heatmap we need to add two for loops:

Seaborn makes it way easier to create a heatmap and add annotations:

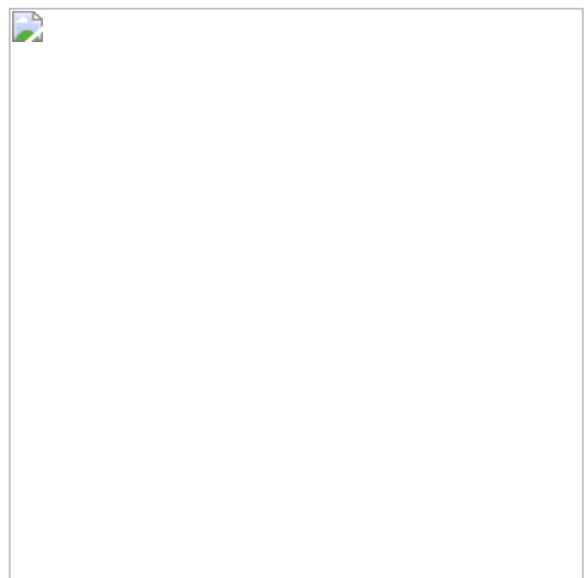
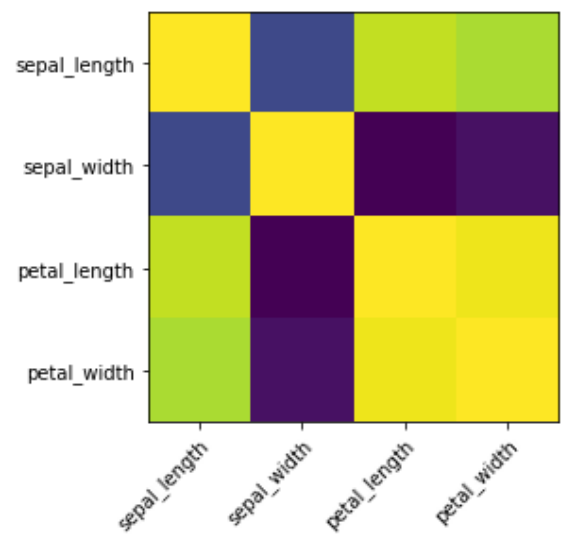


Figure 23: Heatmap without annotations

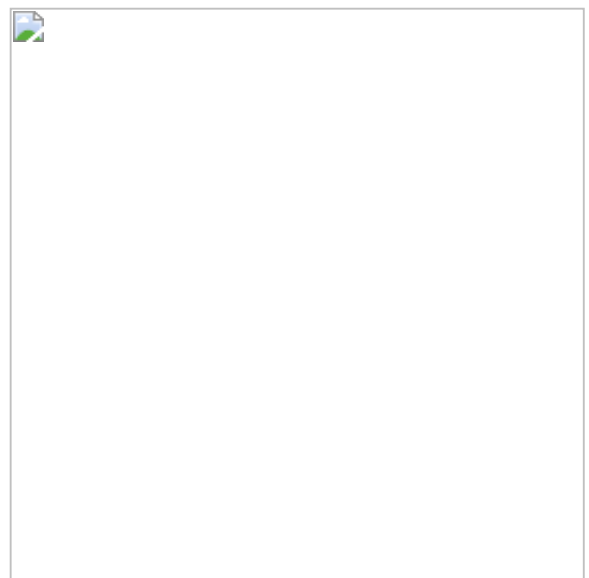
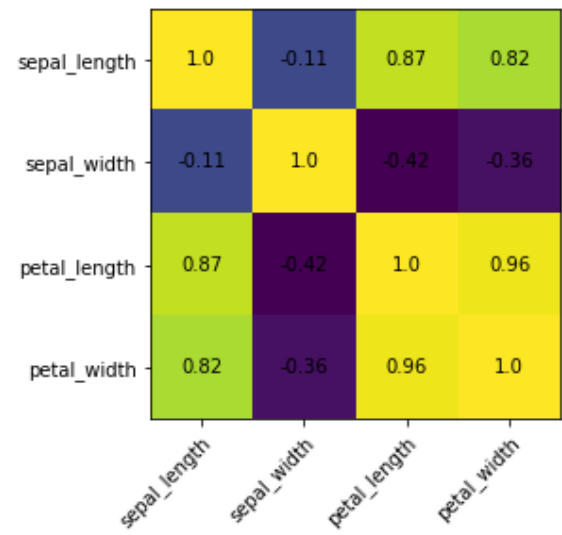


Figure 24: Heatmap with annotations

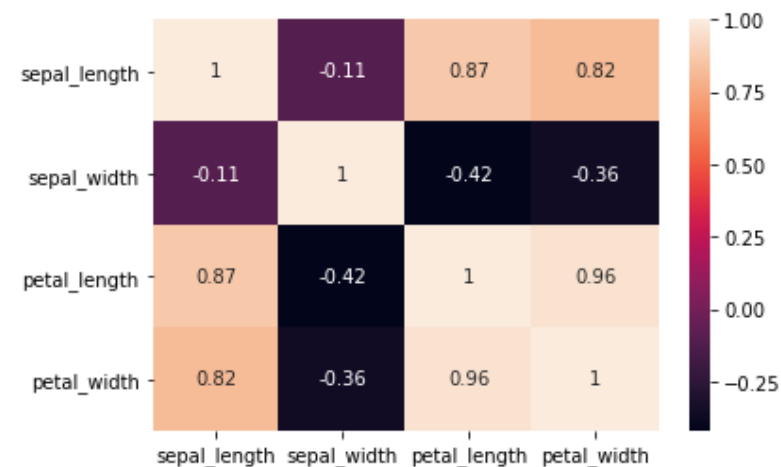


Figure 25: Heatmap with annotations

Faceting

Faceting is the act of breaking data variables up across multiple subplots and combining those subplots into a single figure.

Faceting is really helpful if you want to quickly explore your dataset.

To use one kind of faceting in Seaborn we can use the `FacetGrid`. First of all, we need to define the `FacetGrid` and pass it our data as well as a row or column, which will be used to split the data. Then we need to call the `map` function on our `FacetGrid` object and define the plot type we want to use, as well as the column we want to graph.

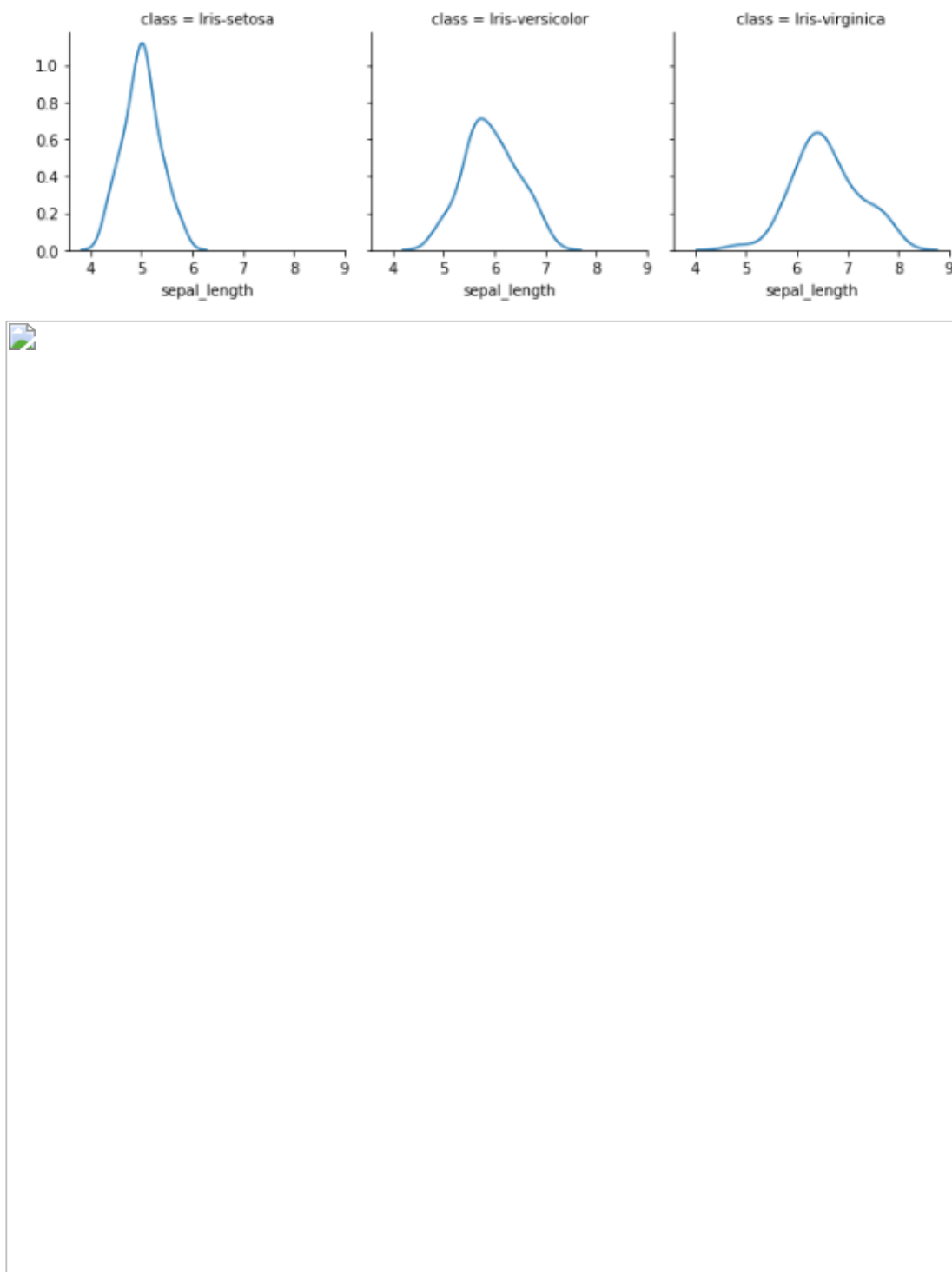
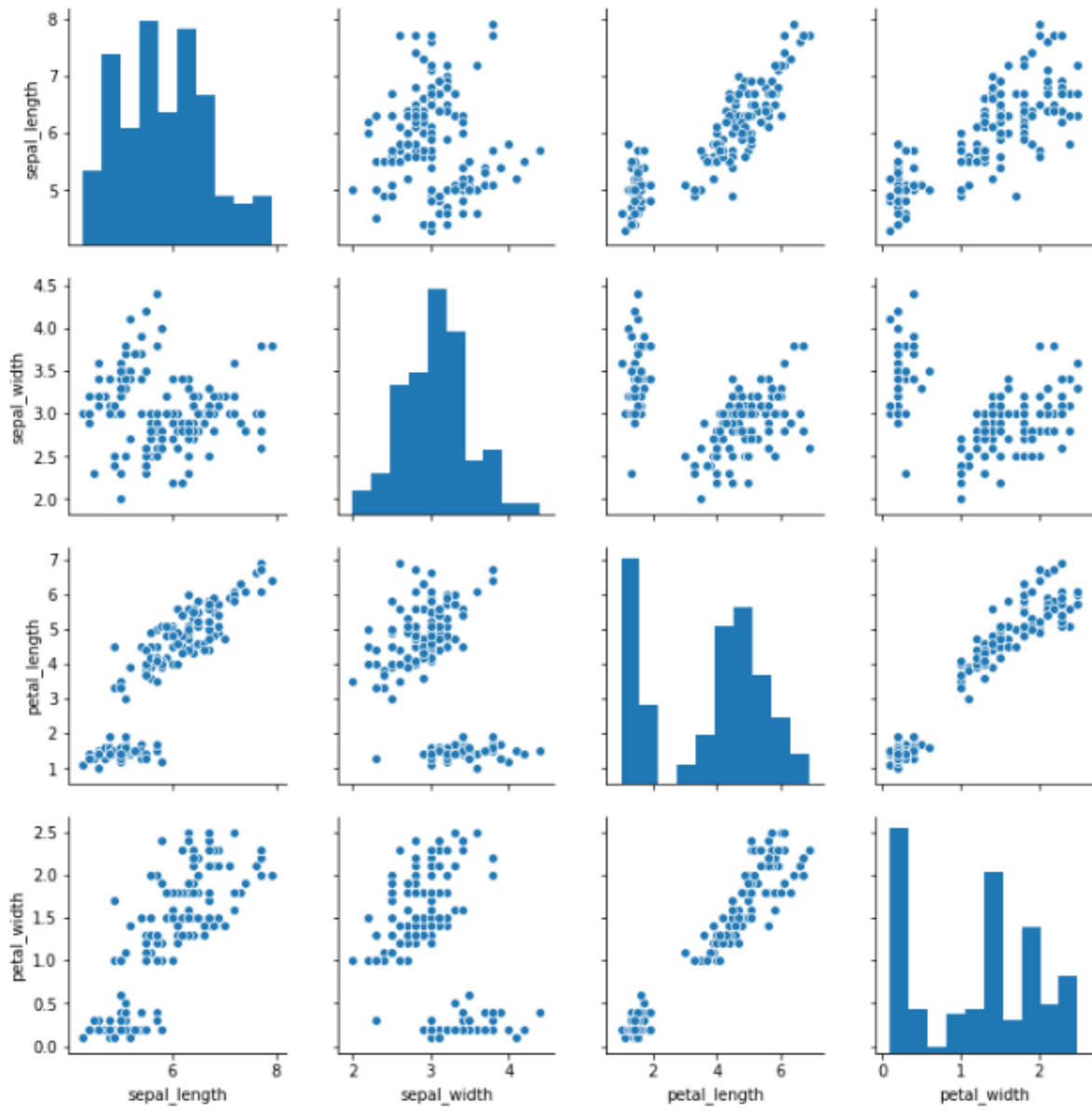


Figure 26: Facet-plot

You can make plots a lot bigger and more complicated than the example above. You can find a few examples [here](#).

Pairplot

Lastly, I will show you Seaborn's `pairplot` and Pandas' `scatter_matrix`, which enable you to plot a grid of pairwise relationships in a dataset.



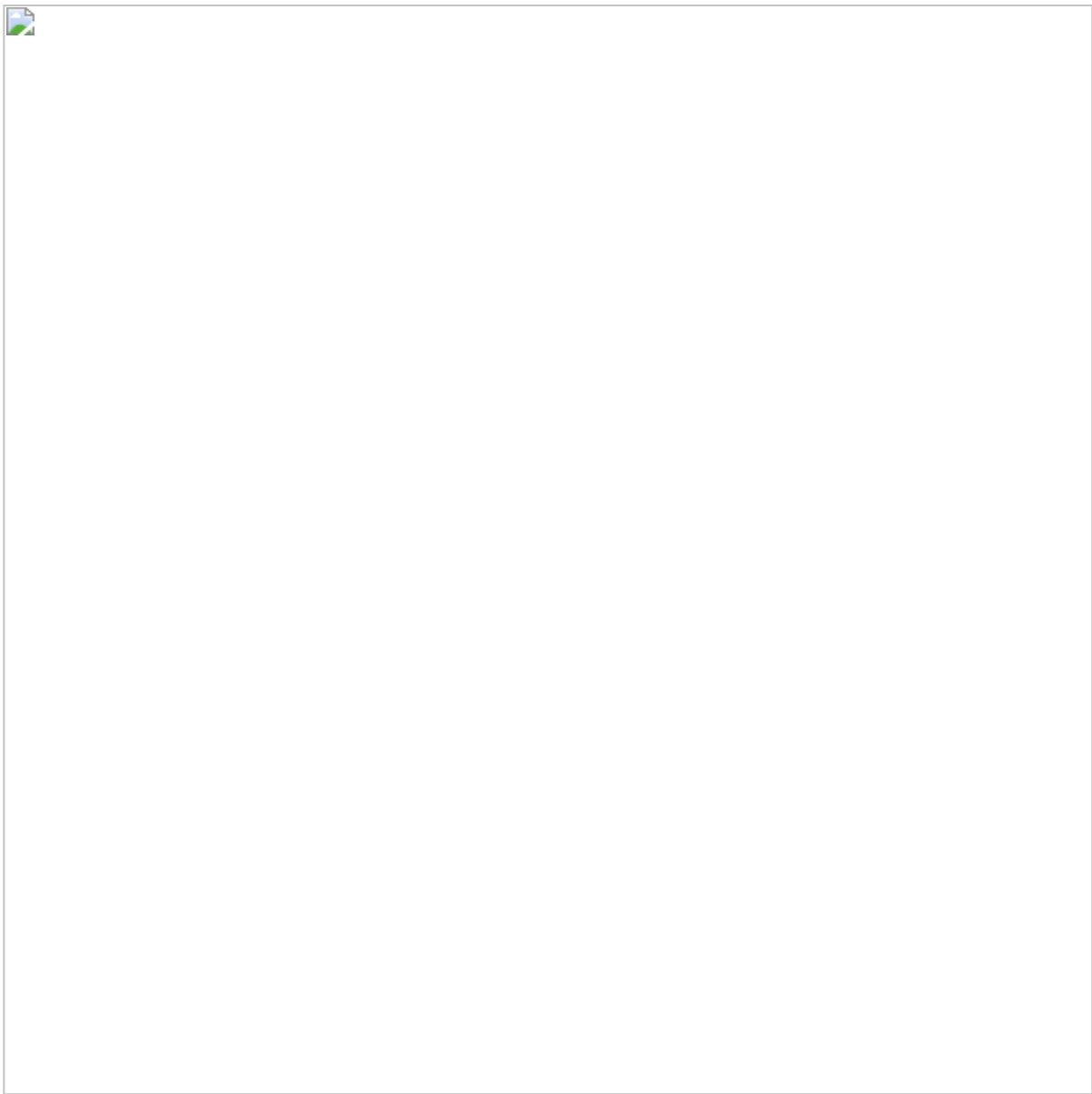
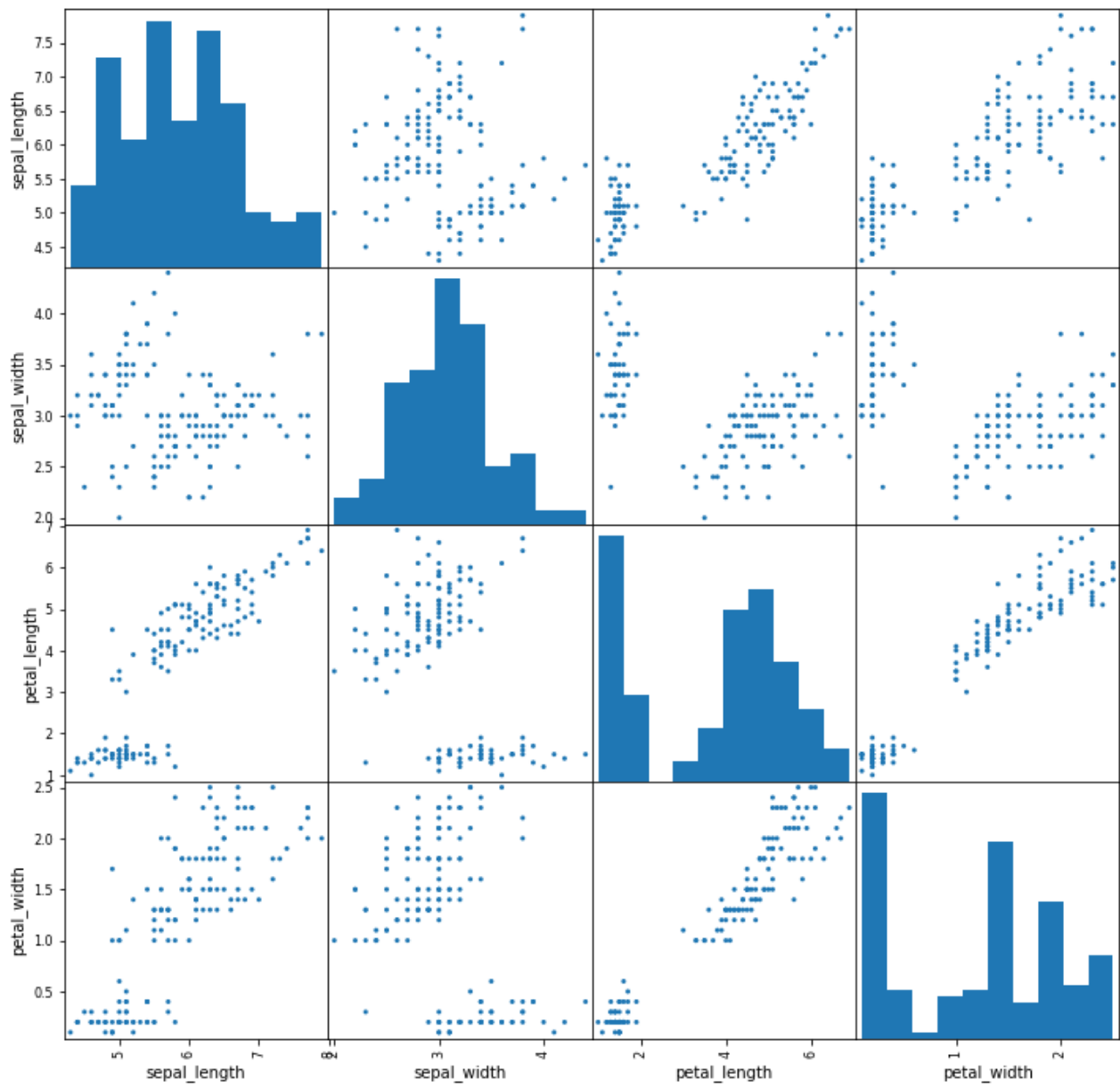


Figure 27: Pairplot



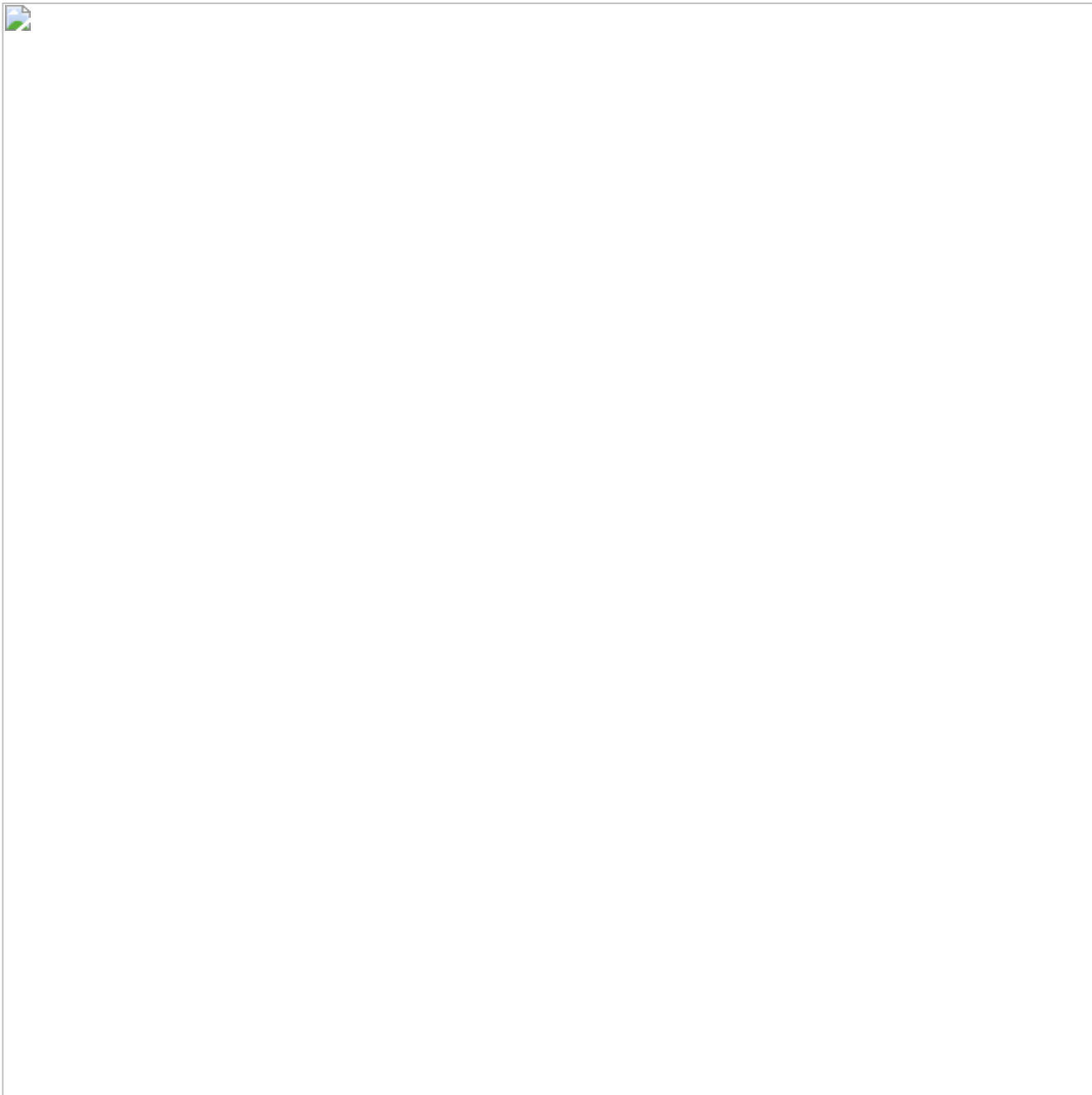


Figure 28: Scatter matrix

As you can see in the images above these techniques are always plotting two features with each other. The diagonal of the graph is filled with histograms and the other plots are scatter plots.

Recommended Reading

Scraping Reddit data

How to scrape data from Reddit using the Python Reddit API Wrapper(PRAW)

towardsdatascience.com

Conclusion

Data visualization is the discipline of trying to understand data by placing it in a visual context so that patterns, trends and correlations that might not otherwise be detected can be exposed.

Python offers multiple great graphing libraries that come packed with lots of different features. In this article, we looked at Matplotlib, Pandas visualization and Seaborn.

If you liked this article consider subscribing on my [Youtube Channel](#) and following me on social media.

The code covered in this article is available as a [Github Repository](#).

If you have any questions, recommendations or critiques, I can be reached via [Twitter](#) or the comment section.